

CLAIMS

What is claimed is:

1. A method for reducing IOs by coalescing writes in a computer system, comprising:
identifying, in a first storage location, a first data block ready to be written into a second storage location, the first data block having a first data block address;
identifying, in the first storage location, additional data blocks to be written into the second storage location; and
writing the identified first and additional data blocks to the second storage location with a single write IO,
wherein the first data block and the additional data blocks form a set of data blocks with consecutive data block addresses to be stored in the second storage location, the consecutive data block addresses having a range that contains the first data block address.
2. The method of claim 1, further comprising tracking a total number of the identified first and additional data blocks.
3. The method of claim 2, further comprising setting a predetermined upper limit for the total number of the identified first and additional data blocks, wherein if the predetermined upper limit is met by the total number of the identified first and additional data blocks, the method stops identifying additional data blocks and immediately writes the identified first and additional data blocks to the second storage location.
4. The method of claim 1, further comprising copying each of the identified first and additional data blocks to a temporary storage location, wherein writing the identified first and additional data blocks to the second storage location is accomplished by writing copies of the identified first and additional data blocks in the temporary storage location to the second storage location.

5. The method of claim 4, further comprising marking the identified first and additional data blocks in the first storage location not dirty each time after an identified data block is copied to the temporary storage location.
6. The method of claim 5 in which the temporary storage location is a temporary buffer cache.
7. The method of claim 1, wherein identifying additional data blocks comprises:
sequentially searching the first storage location to identify next data blocks with higher and lower data block addresses consecutive to the first data block address; and
determining whether an identified next data block is dirty,
wherein if the identified next data block is not dirty, the method stops sequentially searching for new data blocks with data block addresses at a same data block address side to the data block address of the identified next data block.
8. The method of claim 7 in which sequentially searching the first storage location is conducted alternatively at the lower and higher data block address sides to the first data block address.
9. The method of claim 8 in which sequentially searching the first storage location begins by searching the next data block at the lower data block address side.
10. The method of claim 8 in which sequentially searching the first storage location begins by searching the next data block at the higher data block address side.
11. The method of claim 7 in which sequentially searching the first storage location comprises:
searching for all dirty next data blocks with consecutive data block addresses to the first data block address at one of the lower and higher data block address sides; and

searching for all dirty next data blocks with consecutive data block addresses to the first data block address at another one of the lower and higher data block address sides.

12. The method of claim 1 in which the first storage location is a buffer cache.
13. The method of claim 1 in which the second storage location is a disk.
14. The method of claim 1 in which the computer system is a database system.
15. The method of claim 1 further comprising:
copying the first and additional data blocks to form a copy of the first and additional data blocks; and
allowing other entities in the computer system to access the copy of the first and additional data blocks during the act of writing the first and additional data blocks.
16. The method of claim 1 wherein identifying the additional data blocks comprises:
sequentially searching the first storage location to identify next data blocks in the direction of either the higher or lower data block addresses consecutive to the first data block address.
17. The method of claim 16 wherein if an identified next data block is not dirty, sequentially searching in the opposite direction of either the higher or lower data block addresses consecutive to the first data block address.
18. The method of claim 1 in which the computer system comprises a database system.
19. The method of claim 1 in which the first storage location is a buffer cache, wherein the entire buffer cache is searched.

20. The method of claim 19 in which hashing is performed to search the entire buffer cache.
21. A computer system, comprising:
first storage means for storing data in a first plurality of data blocks;
second storage means for storing data in a second plurality of data blocks;
means for identifying, in the first storage means, a first data block and additional data blocks to be written into the second storage means, the first data block having a first data block address;
means for writing the identified first and additional data blocks to the second storage means with a single write IO,
wherein the first data block and the additional data blocks form a set of data blocks with consecutive data block addresses to be stored in the second storage means, the consecutive data block addresses having a range that contains the first data block address.
22. The computer system of claim 21, further comprising:
means for tracking a total number of the identified first and additional data blocks;
and
means for setting a predetermined upper limit of the total number of the identified first and additional data blocks,
wherein if the predetermined upper limit is met by the total number of the tracking means, the computer system stops identifying additional data blocks and immediately writes the identified first and additional data blocks to the second storage means.
23. The computer system of claim 21, further comprising:
temporary storage means for temporarily storing data in a third plurality of data blocks;
means for copying each of the identified first and additional data blocks to the temporary storage means; and

means for marking the identified first and additional data blocks in the first storage means not dirty each time after an identified data block is copied to the temporary storage means,

wherein the identified first and additional data blocks are written to the second storage means by copying copies of the identified first and additional data blocks in the temporary storage means to the second storage means.

24. The computer system of claim 23, wherein the temporary storage means is a temporary buffer cache.

25. The computer system of claim 21, wherein said identifying means comprises:
means for sequentially searching the first storage means to identify the first data block and next data blocks with higher and lower data block addresses consecutive to the first data block address of the first data block; and

means for determining whether the an identified next data block is dirty,
wherein if the identified next data block is not dirty, the computer system stops sequentially searching for new data blocks with data block addresses at a same data block address side to the data block address of the identified next data block.

26. The computer system of claim 25, wherein the sequentially searching means searches the first storage means alternatively at the lower and higher data block address sides to the first data block address.

27. The computer system of claim 25, wherein the sequentially searching means searches the first storage means for all dirty data blocks with consecutive data block addresses to the first data block address at one of the lower and higher data block address side before it searches the first storage means for all dirty data blocks with consecutive data block addresses to the first data block address at another one of the lower and higher data block address sides.

Express Mail Label No. EV348160486US

PATENT
OI7030762001

28. The computer system of claim 21 in which the first storage means is a buffer cache.
29. The computer system of claim 21 in which the second storage means is a disk.
30. The computer system of claim 21 in which the computer system is a database system.
31. A computer program product comprising a computer usable medium having executable code to execute a process for reducing IOs by coalescing writes in a computer system, the process comprising:
- identifying, in a first storage location, a first data block ready to be written into a second storage location, the first data block having a first data block address;
 - identifying, in the first storage location, additional data blocks to be written into the second storage location; and
 - writing the identified first and additional data blocks to the second storage location with a single write IO,
- wherein the first data block and the additional data blocks form a set of data blocks with consecutive data block addresses to be stored in the second storage location, the consecutive data block addresses having a range that contains the first data block address.